



Understanding and Exploiting the Energy advantages of Field-Programmable Technologies

André DeHon andre@seas.upenn.edu



Outline

- Energy: why care?
- Teasers
- Architecture
- Implication for
 - Applications
 - FPGA Architecture
- Information (Coding)
- Variation and Aging

Energy Setup

Energy



- Growing domain of portables
 - Less energy/op \rightarrow longer battery life
- Global Energy Crisis
- Power-envelope at key limit
 - E reduce \rightarrow increase compute in P-envelope
 - Scaling
 - Power density **not** transistors limit sustained ops/s
 - Server rooms
 - Cost-of-ownership not dominated by Silicon

-Cooling, Power bill

Energy, Reliability, Capacity Squeeze

- Smaller devices → higher variation → increased voltage margins

 ITRS shows mostly flat V_{dd}
- Capacitance/gate decreases little
 - Maybe halved 45nm to 11nm [ITRS]
 - When Gate count increase 16× ?!?
- Power density already limiting use of transistor capacity on die

Trends threaten ability to exploit smaller technologies.



Dopants

 $\propto N$



- Based on ITRS2009 data
- Assume start power-limited designs, gap grows with continued scaling $P \propto N_{gates} C_g V^2 f_6$



DeHon--FPT 2015

http://www.nap.edu/catalog.php?record_id=12980

Variation threatens E/Op reduction





Min-Energy for multiplication (typically subthreshold)

[Bol et al., IEEE TR VLSI Sys 17(10):1508—1519]

Energy Limited



- It is Energy that defines
 - Ops/s can extract from a power-limited chip
 - Ops/battery-hour can extract from a portable
- If a technology makes E/op worse
 - That technology is worse

-End-of-scaling



Energy

$E \propto \alpha C V^2$

- C driven by architecture
 - Also impacted by variation, aging
- V today, driven by variation, aging
- α driven by architecture, coding/information

Energy

Architecture Variation/Aging $E \propto \alpha C V^2$ Information

Teasers

Processors and Energy

- Very little into actual computation
- Determine and Fetch Instruction
- Read and Write data from memories





FPGA Energy Advantage

- No memory energy
- Exploit data correlation (Information)
- Performance with
 low clock rate
- Fine-grained specialization
 - Just the needed operations



[Dally et al. / Computer 2008]

Energy-Efficiency of Spatial Computation

- Considerable evidence of lower energy on FPGAs than Processors and GPUs
- SPICE simulation example

- V6LX760 (40nm) vs. Core i7 965 (45nm)



More Preliminary Evidence

- Microsoft Catapult justifies FPGAs in Cloud servers by Energy reduction
- FPGAs for finance driven by limited powerdensity in servers "near" trading floor
- Rush to "accelerators" driven by energy inefficiency of processor
- This conference, FCCM
 Full of papers showing FPGA energy advantage

Energy Advantage FPGAs

- Why?
- Potential for additional advantage?
- How do we increase advantage?

Outline

- Energy: why care?
- Teasers
- Architecture
- Implication for
 - Applications
 - FPGA Architecture
- Information (Coding)
- Variation and Aging



Architecture

E(spatial)<E(sequential)

E(FPGA)<E(processor)

Architecture Outline

- Rent's Rule and VLSI Complexity
- Memories
- Central Processor
- Spatial Locality
 - Data
 - Instruction
- Wire Sharing

Bisection Width

- Partition design into two equal size halves – Minimize wires (nets) with ends in both halves
- Number of wires crossing is bisection width
 - Information crossing
- lower bw = more locality



Rent's Rule

 If we recursively bisect a graph, attempting to minimize the cut size, we typically get:

$BW=IO = c N^p$

- –0≤p≤1
- –p≤1 means many inputs come from within a partition

[Landman and Russo, IEEE TR Computers p1469, 1971]

DeHon--FPT 2015



Rent and Locality

- Rent and IO quantifying locality
 - local consumption
 - local fanout

 $IO = c N^p$



Common Applications

- Rent p=0

 Shift-register, 1D filter
- Rent p=0.5
 - Array multiplier
 - 2D Window Filter
 - nearest-neighbor
- Rent p=1.0 – FFT, Sort



VLSI Interconnect Area

- Bisection width is lower-bound on IC width
 - When wire dominated, may be tight bound
- (recursively)
- Rent's Rule tells us how big our chip must be



Rent Network Richness



Architecture Outline

- Rent's Rule and VLSI Complexity
- Memories
- Central Processor
- Spatial Locality
 - Data
 - Instruction
- Wire Sharing

Memory Energy

- Reading out of large
 memories expensive
- The larger the memory, the more expensive per bit read
- O(M^{0.5}) wire length
- Random access must send address

 O(M^{0.5}log(M))
- Sequential access
 - $O(M^{0.5})$ per bit



Processors and Energy

- Most energy reading out of two memories
 - Instruction

Data





Central Processor with **Description Locality**

- p<1.0, total instruction
 Read/write O(N) bits bits are O(N)
 - $O(N^{1.5}log(N))$



Instruction Sharing

- **Problem:** absolutely, instruction energy dominates data energy
- **Assumption:** every bit operator needs its own unique instruction
- **Opportunity:** share instructions across operators

- Looping, wide words

• Still O(N^{1.5}) data memory energy

Instruction Sharing (I) (p=0.7)

Energy Ratio to Processor W=1, I=N



SIMD+Instr. Share (p=0.7)

Energy Ratio to Processor W=1, I=N



Architecture Outline

- Rent's Rule and VLSI Complexity
- Memories
- Central Processor
- Spatial Locality
 - Data
 - Instruction
- Wire Sharing

Data Locality

- **Problem:** Must pay $O(N^{0.5})$ for every read since data must be moved in and out of memory. \sqrt{M}
- **Opportunity:** compute local to data



Sequential with Data Locality

- Place for locality --Rent Partitions
- Store data at endpoints
- Send through network from producer to consumer
- Store location at leaves – O(log(N))
- Build H-Tree to keep area to O(Nlog(N))


Sequential with Data Locality

- Area = O(Nlog(N))
- Sending addresses log(N) bits at top
- Signals lower O(1)
- Only send a few over top O(N^p)
- O((log^{1.5}N)N^{p+0.5}) for p>0.5
- Cheaper to send where needed than to central location.



Data Local Compare (p=0.7)



Data Local Compare (p=0.7)



Data Local Compare (p=0.7)



Instruction Locality

- Problem: Multiply energy by O(log(N)) to send an address up the tree
- Opportunity: store instructions local to switches
 - In tree
- ...what an FPGA does!



41

Fully Spatial (FPGA)

- An FPGA
- Each signal gets own wire
- No addresses
- Configuration local
- Area grows as O(N^{2p}) for p>0.5
- Energy O(N^{2p}) for p>0.5
 - Θ(N) for p<0.5







p=0.7 Compare to FPGA

Energy Processor/FPGA



DeHon--FPT 2015

Asymptotic Energy

Org	Any p	p<1.0	1>p>0.5	p=0.5	p<0.5		
Processor	O(N ^{1.5} log ^{1.5} N)	O(N ^{1.5} logN) Description Locality					
FPGA 2-metal		(D(N ^{2p})	O(Nlog ² N)	Θ(N)		

Note break at p=0.75

p=0.8 Compare to FPGA

Energy Processor/FPGA



Compare around p=0.75

• Breakpoint at p=0.75



Intuition

- Given a good spatial layout
 - It is cheaper to transmit the result of a gate to its well-placed consumers
 - Fixed metal layers: average wire length o p<0.5: O(1) $\circ p < 0.75$: $O(N^{2p-1}) < O(N^{0.5})$ \circ p>0.75: O(N^{2p-1}) > O(N^{0.5})
 - Than to
 - Fetch inputs from a large central memory ○ O(N^{0.5})





Architecture Outline

- Rent's Rule and VLSI Complexity
- Memories
- Central Processor
- Spatial Locality
 - Data
 - Instruction
- Wire Sharing

Wiring Dominates

Problem: When p>0.5 wiring dominates area

 \rightarrow force longer wires

• **Opportunity:** Share wires

But also, keep instructions local to switches

Instructions Local to Switches

- Constant metal
- Build p<0.5 tree
- Store bits local to each tree level
- Read out of memory there
- Bits/switch differs
 with tree level
- Signal on wire dominates reads
- O(N^{p+0.5}) for p>0.5









Results: Energy

Org	Any p	p<1.0	1>p>0.5	p=0.5	p<0.5		
Processor	O(N ^{1.5} log ^{1.5} N)	O(N ^{1.5} IOgN) Description Locality					
Data Locality (Packet Switch)		0(N	I ^{p+0.5} Iog ^{1.5} N)	O(Nlog ^{2.5} N)	O(Nlog ^{1.5} N)		
FPGA 2-metal			O(N ^{2p})	O(Nlog ² N)	Θ(N)		
Multicontext		C)(N ^{p+0.5})	O(NlogN)	Θ(N)		





Architecture

• There are architectural energy advantages to FPGAs over processors





Optimizing FPGAs and Applications

Tune Parallelism

- For many, regular designs we can tune the parallelism in the implementation on top of the FPGA
- Parameterize number of PEs
- Store data in Embedded Memories
- Tune to balance Memory vs.
 Interconnect costs

Optimal Level of Parallelism

- For regular designs
 - Build common Operator
 - Share amongst logical operations
 - Use embedded RAMs to hold state of multiple operators



62

FPGA Energy Efficiency

- When data memory dominates
- Question is not:
 - How efficient is FPGA vs. ASIC? ... but
 - How efficient is FPGA memory vs. ASIC memory?



FPGA Memory Energy Optimization

- Energy inefficiency from
 - Memory blocks too large
 - Memory blocks too frequent, infrequent
 - Activating memory unnecessarily
- Can get within factor of 2
 - Appropriate distribution of blocks
 - Efficient memory banking

Modern FPGA



- Memory Bank size? (Distribution of sizes?)
- Memory Bank frequency?

Matrix Multiply, Single Memory

Energy normalized to best match



[Kadric et al./FPGA2015] 66

Internal Banking

- Recall memory energy scale as O(M^{0.5})
- Continuous Hierarchy Memory

 Internal Banking
- Only pay as much energy as need



[Kadric et al./FCCM2014+FPGA2015]



Information (Switching Activity)



Encode to Reduce α

- Dominant FPGA energy
 - Is energy in wires

Activity Reduction (%)

40

30

20

10 -

0

- especially in p>0.5 cases
- Encode to reduce switching



stereovision1

DeHon--FPT 2015

Outline

- Energy: why care?
- Teasers
- Architecture
- Implication for
 - Applications
 - FPGA Architecture
- Information (Coding)
- Variation and Aging

 $E \propto \alpha C V^2$

Variation and Aging
Voltage Challenge

- End of Dennard Scaling
 - Subthreshold slope prevents linear reduction of Voltage with feature size
- Variation
 - Devices no longer "identical"
 - Typically accommodate with higher voltage
- Aging
 - Devices change
 - Typically accommodate with higher voltage

Inflection Point Collision Defeat Scaling

- Spend energy for reliability
 - Margins
 - Replication
- Reduce or eliminate net benefits of feature size reduction



Variation threatens E/Op reduction





[Bol et al., IEEE TR VLSI Sys 17(10):1508—1519]



- High margins driven by uncommon tails
 Most devices much better
- Large device count → sample further
 into tails



PDF

 V_{th}

Mehta 2012: 22nm PTM LP – 10,000 samples 76

DeHon--FPT 2015

Wear Out

- Reduced burnin opportunity
- Aging: NTBI, HotCarrier, Electromigration, ...
- Forces margin for End-of-Life
 - ...and that may not be enough



Post Fabrication Configurability

 ASICs bind functions to physical transistors **before** fabrication

Before know behavior of device

 FPGAs bind functions to physical devices after fabrication

After device characteristics determined

Variation Tolerance

- Idea: assign resources, post fabrication to compensate for variations
- **Opportunity**:
 - Balance fast paths and slow paths
 - Assign slow resources to non-critical paths
 - Avoid devices in uncommon tails
 - Scale voltage down more aggressively
- Fixed design limited to worst-case path

 Must scale voltage up so path meets timing
- Paradigm shift: Component-specific mapping
 DeHon--FPT 2015

Variation Challenge



- Use of high V_{th} resource forces high supply voltage (V_{dd}) to meet timing requirement
- Delay: CV/I and I goes as $(V_{dd}-V_{th})^2$

Component-Specific



- Avoid high V_{th} resource
- Allow lower supply voltage (V_{dd}) to meet timing requirement
- Delay: CV/I and I goes as $(V_{dd}-V_{th})^2$

Component-Specific Assignment



- Could come out other way
 - Best mapping unique to component



Lifetime Adaptation

- Issues: Devices will wear out during operation
 - Parameters vary with time
 - *E.g.* resistance increases \rightarrow slower operation
- **Opportunity:** reconfigure to replace bad or slow devices
- **Detect:** Concurrent error detection, RAZOR, lightweight application checks, periodic self test
- **Recover:** reassign resources

- re-invoke mapping

• **Paradigm shift:** mapping throughout lifetime

Energy vs V_{dd} (des)



[Mehta, FPGA 2012]

Energy vs V_{dd} (des)



[Mehta, FPGA 2012]

Energy vs V_{dd} (des)

- Nominal uses minimum size
- Delay-aware routing reduces energy margins
 - 1. Smaller sizes



[Mehta, FPGA 2012]

Energy vs V_{dd} (des)

- Nominal uses minimum size
- Delay-aware routing reduces energy margins
 - 1. Smaller sizes
 - 2. Lower voltages



[Mehta, FPGA 2012]

87

Energy vs V_{dd} (des)

- Nominal uses minimum size
- Delay-aware routing reduces energy margins
 - 1. Smaller sizes
 - 2. Lower voltages
 - 3. Less Leakage



[Mehta, FPGA 2012]

88

Minimum Energy vs Technology

- Delay-oblivious scales to 16nm
- Delay-aware scales to 12nm at least....
- Extend useful life of Silicon technology generation



DeHon--FPT 2015

[Mehta PhD thesis 2012 (refined from FPGA 2012)] 89

Delay Map



• LAB (27,22)

[Gojman, FPGA2013] DeHon--FPT 2015



DUK Delay (ps)

Choose Your own Adventure

- Idea: Precompute Alternate mappings
 - Still just one bitstream



[Rubin, FPGA 2008] 92



DeHon--FPT 2015

Defect-Level Viable with FPGAs



- Fine-grained repair
- Avoiding routing defects
 - Tolerates >20%
 switch defects
 - Ongoing work
 - Different defect types and resources
- Don't have to sample tails



[Rubin/unpublished]

93

Lifetime Failure?

Go back to alternatives and load again.

Maybe just the one that failed.



[Giesen et al./unpublished] 94

– Millisecond repair.

DeHon--FPT 2015

$E \propto \alpha C V^2$

Wrapup

Fund. Energy Benefits of FPGAs

- Architecture
 - Minimize energy in data movement
 - Asymptotic advantage
 - Bound effects of mismatch
- Application



- Tune parallelism to problem and problem size
- Specialize to application
- Information exploit correlation, switching
- Post-fabrication configuration
 - Avoid tails of variation and aging

– Smaller features, lower voltage





ic.ese.upenn.edu

Papers: Proc. IEEE 2015: Fundamental Underpinnings TRETS 2014: GROKLAB, FCCM2014: GROKINT FCCM2014, FPGA2015: Kung Fu, Mem. Arch. FPGA 2012: Component-Specific Routing **TRETS 2011: CYA 211** 97

DeHon--FPT 2015